

Interaktionen und Modifikation von RNA und Proteinen

Prof. Sonja Prohaska

Computational EvoDevo
Universität Leipzig

SS 2018

DNA is “never” naked in a cell

general binder

DNA is usually in association with proteins. In all domains of life there are small, basic **chromatin associated proteins (ChAP)** attached to the DNA forming higher order structures. Their function is to prevent DNA from agglutination, ensuring stability and flexibility, aid structure formation (“packaging”) and engage in gene regulation.

specific binder

Sequence-specific **transcription factors** associate with specific binding sites. Their functions is commonly gene regulation.

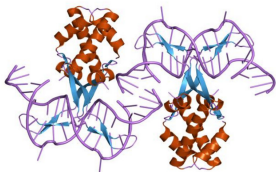
Proteins that are general binders of DNA

eukaryots

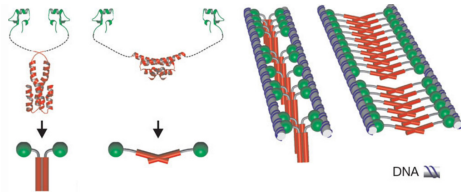
- ▶ the DNA-protein complex: chromatin
- ▶ major chromatin associating proteins (ChAP): **histones**

bacteria

- ▶ the DNA-protein complex: nucleoid
- ▶ nucleoid associating proteins (NAP):
 - HU** – introduces negative supercoiling
 - H-NS** – regulates gene expression



HU, core (red), arms (blue)



histone-like nucleoid-structuring protein (H-NS)

general vs. specific binder

general binder

- ▶ small basic proteins
- ▶ contact negatively charged DNA backbone
- ▶ e.g. HU, Alba, histones (double stranded)
- ▶ e.g. SSB (single stranded)
- ▶ bind everywhere, often bend the DNA

specific binder

- ▶ large protein with DNA-binding domain(s) (DBD)
- ▶ contact major groove of DNA
- ▶ DBDs: e.g. HTH (helix-turn-helix), zinc finger, leucine zipper, helix-loop-helix
- ▶ bind to specific sites
- ▶ sequence-dependent

DNA-binding proteins – a molecular view

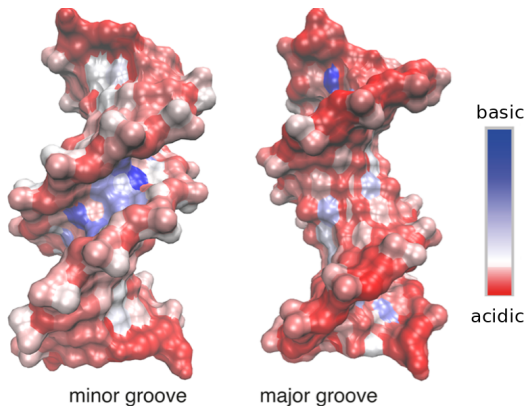
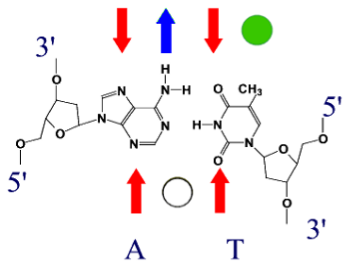
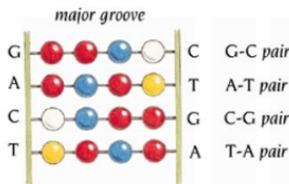


Fig: Electrostatic potential surface of the DNA. General binders are basic and contact the acidic DNA backbone. Specific binders make contact with the nucleobases in the major and minor groove.

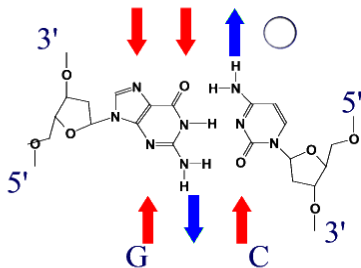
Major Groove



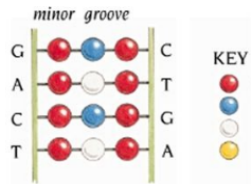
Minor Groove



Major Groove



Minor Groove



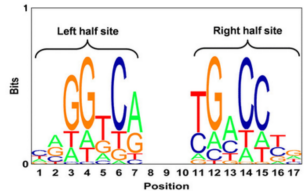
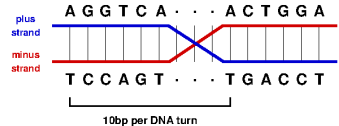
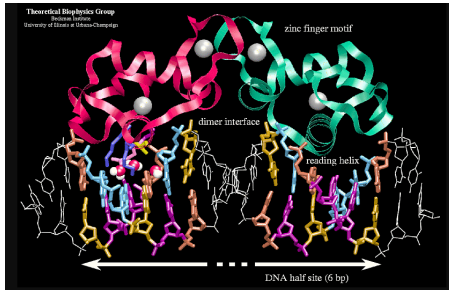
KEY

- = H-bond acceptor
- = H-bond donor
- = hydrogen atom
- = methyl group

Proteins can distinguish all four nucleobases in the major groove. In the minor groove only AT and GC basepairs can be distinguished.

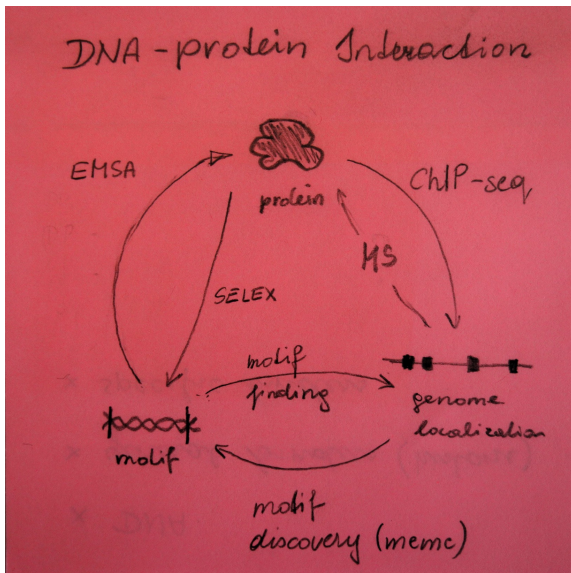
DNA-binding proteins – a molecular view

Example: the estrogen receptor dimer



Homodimeric proteins consist of two identical units that are in contact via the dimer interface. Here each dimer makes contact to 6 base pairs in the major groove in a sequence-specific manner. The motif is **palindromic**.

How to study DNA-protein binding?



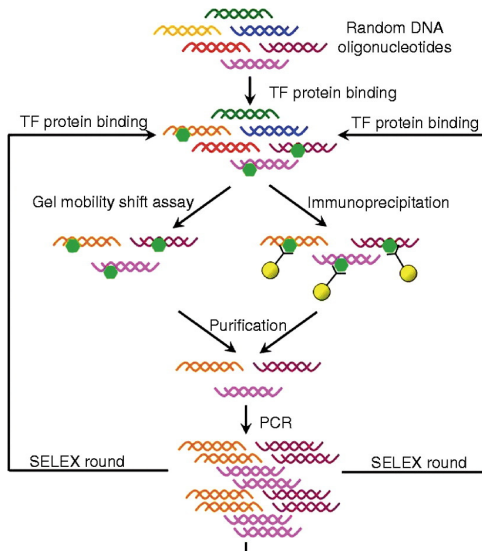
How to study DNA-protein binding?

Given a protein, which DNA sequence does it bind preferentially?

in vitro selection – SELEX

- ▶ SELEX (systematic evolution of ligands by exponential enrichment) also known as in-vitro evolution
- ▶ a very large oligonucleotide library
 - ▶ randomly generated sequences of fixed length
- ▶ is exposed to the target protein
- ▶ unbound oligos flow through (are removed)
- ▶ bound sequences are eluted and amplified by PCR
 - ▶ another SELEX round with increased stringency is done or
- ▶ the pool of sequences is finally sequenced
- ▶ an alignment reveals the **motif** and its variation

How to study DNA-protein binding?



	BanII		PstI
	GGATCCCCGCAGG---
1,	-----G	TGAGTCA	CCATCCCGGTTGGC---
2,	-----GACAGGC	TGAGTCA	TATGCACG-----
3,	-----GAGCAA	TGAGTCA	TCGTGTCCGG-----
4,	-----GA	TGAGTCA	CAGGCAACGGGCAC---
5,	-----A	TGACTCA	TTGAGCGACTTACCG---
7,	CTGTATTACGACACAA	TGAGTCA	
8,	-----ATAACGAGTGGGA	TGACTCA	TT-----
9,	-----ATTCTCGCCTTCTG	TGAGTCA	T-----
11,	-----AAAAAA	TGACTCA	TCCGAGCT-----
12,	-----AT	TGACTCA	TACGCAGCTAGACT---
13,	-----GAGGTAGA	TGACTCA	TGGACTG-----
14,	-----A	TGACTCA	TCGCTAGTCCATGGG---
15,	-----GTGTCCTTCGGGA	TGACTCA	TGC-----
16,	-----GTCACGGGGCC	TGACTCA	TAGAA-----
18,	-----G	TGACTCA	CGGAAATGTTGG---
19,	-----GCATTATGGAG	TGACTCA	TCCT-----
20,	-----GAAGCATTG	TGACTCA	TCGCT-----
21,	-----TCGGTAGTCGGTA	TGAGTCA	TT-----
22,	-----CGG	TGACTCA	CGTAGAGGTAACC---
25,	-----A	TTAGTCA	TCAGGGGTGGCGAC---
26,	-----AAATTTACATCCGA	TGACTCA	TA-----
27,	-----AGCCCGGTGAGAA	TGAGTCA	TA-----
28,	-----TGGCCCCGGGTCTC	TGACTCA	GC-----
29,	-----AACGAGACGCGC	TGAGTCA	TCCT-----
30,	-----GCCAGTTGATATCT	TGTGTCA	CGG-----
32,	-----GTAGCTGAGAG	TGACTCA	CGCCTG-----
33,	-----GGG	TGACTCA	TAAGATAAAATCT---
34,	-----CCCGGTAGGCTCGA	TGACTCA	A-----
35,	-----TCAGATCAGCTTA	TGACTCA	GG-----
36,	-----CGCTGG	TGACTCA	TCGTGTCTG-----
37,	-----GCAGGCGCCACCGG	TGACTCA	T-----
38,	-----GTCAGTTG	TGACTCA	CTCTACC-----
39,	-----GTATGACGTGGA	TGTGTCA	GAGG-----
40,	-----TGGCCCTA	TGACTCA	TAAGGCAC-----
41,	-----GCCAATGACTTCTC	TGAGTCA	T-----
42,	-----CCGTATCGCGGGT	TGAGTCA	CT-----
43,	-----GA	TGAGTCA	GGACCGGGTTGGA---
44,	-----TACGTG	TGACTCA	TTCCACCACCC-----
45,	-----GCCTG	TGACTCA	TCTCGCGGTA-----
46,	-----GAAATA	TGACTCA	CGGACGCTCT-----
47,	-----AGTGTGTAAGGGGG	TGACTCA	T-----
48,	-----AG	TGACTCA	CCGACACGTACTG---
49,	-----TGCGTCGGG	TGACTCA	GATTTG-----

FIG. 3. GCN4-binding sites.

How to represent a motif?

Given an alignment of n motifs of length l
here the result of the SELEX experiment: $n = 43$, $l = 7$

- ▶ **consensus string**

most frequent nucleotide per position 5'-TGAGTCA-3'

- ▶ **consensus pattern**

all possible nucleotides per position 5'-T(G|T)(A|T)(G|C)TCA-3'

$K = \{G, T\}$, $W = \{A, T\}$, $S = \{G, C\} \rightarrow 5\text{'-TKWSTCA-3'}$

- ▶ **position frequency matrix – PFM**

Alphabet = $\{A, C, G, T\}$, size of Alphabet a , matrix $a \times l$

- ▶ **sequence logo** based on the information content (IC)



How to calculate the information content?

The information content (y-axis) is measured in bits.

The maximal information depends on the size k of the Alphabet \mathfrak{A} , here $\mathfrak{A} = \{A, C, G, T\}$, $a \in \mathfrak{A}$, and $k = 4$.

The **Shannon entropy** (uncertainty) H_i at position i of the motif is

$$H_i = - \sum_{\mathfrak{A}} f_i(a) \log_2 f_i(a) \quad (1)$$

where $f_i(a)$ is the relative frequency of nucleotide a at position i .

The **information content** R_i at position i of the motif is

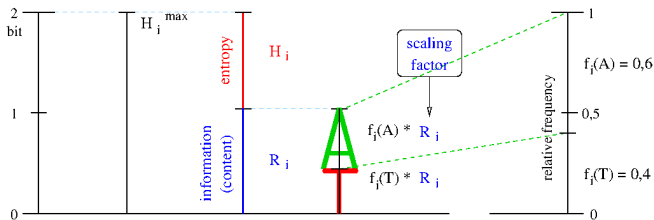
$$R_i = \log_2(k) - H_i \quad (2)$$

(i.e. the maximal information content minus the uncertainty)

The height $h_i(a)$ of a letter a in column i is

$$h_i = f_i(a) \times R_i \quad (3)$$

note: $\log_2 k = \log_2 4 = 2$



... with small-sample size correction

small-sample correction e_n is given by

$$e_n = \frac{1}{\ln 2} \times \frac{s - 1}{2n} \quad (4)$$

where s is the number of sequences in the alignment.

The information content R_i at position i of the motif is then

$$R_i = \log_2(k) - (H_i + e_n) \quad (5)$$

Explain, what does e_n do?

What to do with the motif?

How often would you expect to find a motif m in a sequence M ?

- ▶ assume that motif m is a string, e.g. 5'-GGCCT-3'
 - ▶ of motif length $l = 5\text{bp}$
 - ▶ the sequence M , e.g. the human HoxA cluster,
 - ▶ has a length of $L = 163001\text{bp}$
 - ▶ Alphabet $\mathfrak{A} = \{A, C, G, T\}$
1. assume **uniform** nucleotide distribution: $f(x) = 0.25$
 2. use **mono**-nucleotide distribution:
 $f(A) = 0.2428, f(C) = 0.2555, f(G) = 0.2552, f(T) = 0.2466$
 3. use **di**-nucleotide distribution:
 $f(AG) = 0.0604, f(GG) = 0.0812, f(GC) = 0.0677$
 $f(CC) = 0.0815, f(CT) = 0.0751$

Calculations

$$E^{UNI}(m) = f(X)^5 \times (L - (I - 1)) \quad (6)$$

$$E^{MONO}(m) = (f(G)^2 \times f(C)^2) \times f(T) \times (L - (I - 1)) \quad (7)$$

$$E^{DI}(m) = \frac{f(GG) \times f(GC) \times f(CC) \times f(CT)}{f(G) \times f(C)^2} \times (L - (I - 1)) \quad (8)$$

Results: $E^{UNI}(m) = 159$; $E^{MONO}(m) = 171$; $E^{DI}(m) = 329$;

Which expectation comes closer to the observation?

Let's assume you apply a word search method to the human hox cluster sequence and find 312 motif sites.

Observed sites: number of matches of motif m in M is $O = 312$

Expected sites: $E^{UNI}(m) = 159$; $E^{MONO}(m) = 171$; $E^{DI}(m) = 329$;

Use the **chi-squared test** (χ^2) to determine whether there is a significant difference between the expected frequency and the observed frequency.

$$\chi^2 = \frac{(O - E)^2}{E}$$

df \ p	0.995	0.975	0.9	0.5	0.1	0.05	0.025	0.01	0.005
1	.000	.000	0.016	0.455	2.706	3.841	5.024	6.635	7.879
2	0.010	0.051	0.211	1.386	4.605	5.991	7.378	9.210	10.597
3	0.072	0.216	0.584	2.366	6.251	7.815	9.348	11.345	12.838
4	0.207	0.484	1.064	3.357	7.779	9.488	11.143	13.277	14.860
5	0.412	0.831	1.610	4.351	9.236	11.070	12.832	15.086	16.750

Results: $\chi^2(UNI) = 147.226$; $\chi^2(MONO) = 116.263$; $\chi^2(DI) = 0.878$;

Which expectation comes closer to the observation?

Significance levels

for **one** degree of freedom (df)

expected and observed value, E and O , are significantly different

at 5% level ($p = 0.05$) if $\chi^2 \geq 3.841$

at 1% level ($p = 0.01$) if $\chi^2 \geq 6.635$

Conclusion

$E^{UNI}(m)$ and $E^{MONO}(m)$ are significantly different from O ,

$E^{DI}(m) = 329$ is not.

Therefore we can conclude that the calculation of $E(m)$ using the di-nucleotide distribution is a good prediction for the **observable** number of sites.

DNA is more than a (plus) strand!

- ▶ DNA is double stranded
- ▶ the two strands are referred to as **plus and minus strand**
- ▶ the **convention** is:
 - ▶ the strand which runs from 5' to 3' from left to right
 - ▶ is the top strand (when both strands are displayed) and
 - ▶ it is also the plus strand (the one displayed if only one strand is displayed)
- ▶ the sequence of one strand determines the sequence of the other
- ▶ the relation between the two strands is the “reverse complement”

about motifs and sites

- ▶ the (binding) motif specifies a (binding) pattern
- ▶ a sequence in the genome matching the motif is a (binding) site

Exercises

- ▶ Explain why the expected number of sites E comes closer to the true number of sites O in the sequence M when using the dinucleotide instead of mononucleotide composition.
- ▶ How would you need to adapt equations 6 to 8 to account for both DNA strands in the calculation of the expectation value? What if the motif is palindromic, e.g. 'GGTACC'?
- ▶ How many different motifs of length 5 can be derived from the sequence below?
- ▶ How many binding sites can be found with motif 'GGTCA'?

